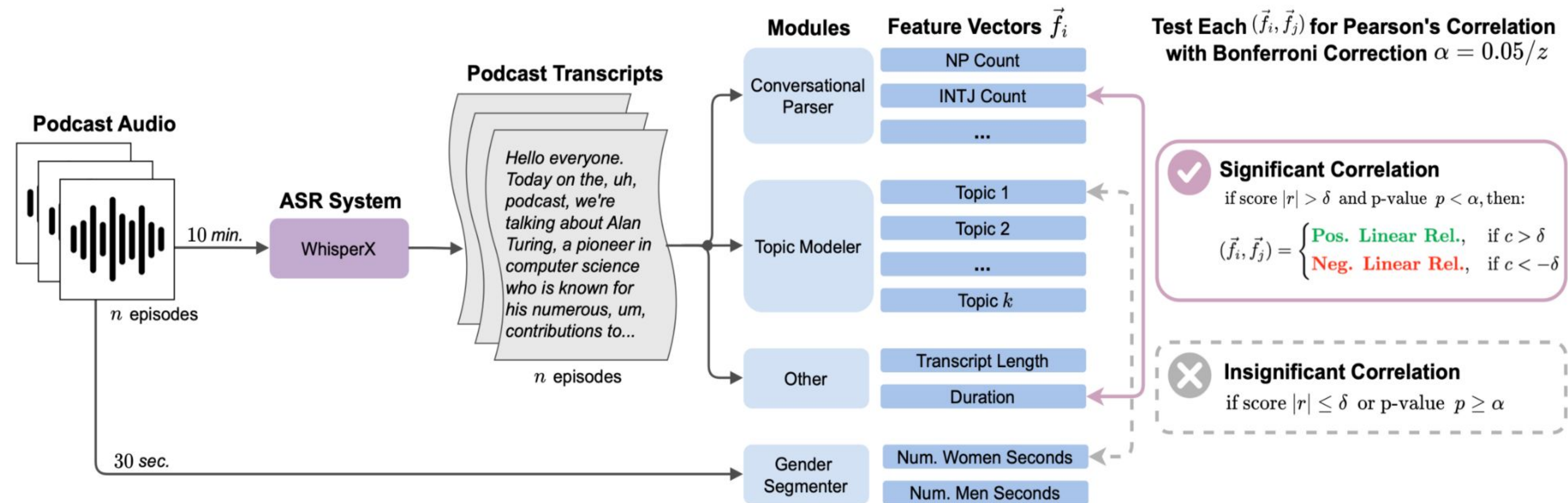# Masculine Defaults via Gendered Discourse in Podcasts and Large Language Models

*Maria Teleki, Xiangjue Dong, Haoran Liu, James Caverlee*

## Women and men's discourse are different.

**(1)** We test for significant correlations between features, including gender and discourse topics.



| Topic N | Gender | $r$ | Topic N Word List | Topic N Categories | Topic N Gender |
|---|---|---|---|---|---|
| Topic 3 | Women | 0.15 | women, woman, men, baby, pregnant, girls, men, doctor, health, birth | Content - Pregnancy | Women |
| | Men | -0.14 | | | |
| Topic 10 | Women | 0.10 | energy, body, feel, mind, space, yoga, love, beautiful, feeling, meditation | Content - Yoga | Women |
| | Men | -0.12 | | | |
| Topic 49 | Women | -0.21 | game, know, think, team, going, mean, play, year, one, good | Content - Sports | Men |
| | Men | 0.17 | | | |
| Topic 71 | Women | 0.14 | christmas, sex, girl, hair, love, get, date, girls, let, wear | Content - Dating | Women |
| | Men | -0.14 | | | |
| Topic 54 | Women | – | get, like, know, right, people, going, podcast, make, want, one | Discourse | Men |
| | Men | 0.12 | | | |
| Topic 60 | Women | -0.27 | going, know, think, get, got, one, really, good, well, yeah | Discourse | Men |
| | Men | 0.20 | | | |
| Topic 62 | Women | 0.33 | like, know, really, going, people, want, think, get, things, life | Discourse | Women |
| | Men | -0.28 | | | |

**(2)** We find that there are discourse topics that have correlations with women or men.

**(3)** This means that, for example, women and men might use a different filler word.

**Men:** $s = And\ I\ was\ \textbf{going},\ hey,\ it's\ cold\ outside...$

**Women:** $s' = And\ I\ was\ \textbf{like},\ hey,\ it's\ cold\ outside...$

## These discourse-based masculine defaults are present in LLM embeddings.

**(4)** So we experiment with flipping the gendered discourse words for (s,s') pairs. We measure the movement of s → s' in the embedding space.

**(5)** We find that men have a more stable/robust embedding representation than women w.r.t. discourse words – this is a representational harm & a masculine default.



$\gamma$: the # of discourse words flipped in s → s'
Percent: Avg. % of S segments which move closer to A_{m,w} after s → s'